



ENVRI Common Operations of
Environmental Research Infrastructures

OEILM: an ontological framework for environmental big data research infrastructures

Zhiming Zhao
z.zhao@uva.nl



**System and Network Engineering,
University of Amsterdam (UvA)**



Project number: 283465

Global warming

Earthquakes

Climate
change

Fresh water

Volcanoes

Deforestation

Epidemic
diseases

Biodiversity loss

Food supplies

Pollution



ENVRI Environmental data and research infrastructure

• Environmental data:

- Observation and measurement, time and locations
- carbon
- plate
- seafloor
- Space weather,
- ocean
- bio diversity
- ...





ENVRI research infrastructure

• Environmental data:

- Observation and measurement, time and locations
- carbon
- plate
- seafloor

Difficulties in sharing and integrating data among different research infrastructures.

• ...

- Environmental research infrastructures
- ENVRI research project
- Open Environmental Information Linking Model (OEILM)
- Use cases
- Discussions

Environmental research ENVRI infrastructure

• An environment research infrastructure

- **Acquisition** -- brings the measures/data streams into the system (**non-reproducible** data)
- **Curation** -- manages/maintains quality data (**reproducible** data)
- **Access** -- facilitates discovery, access (**published** data)
- **Processing** -- facilitates analysis/mining/experiments (**combined/derived** data)
- **Community Support** -- supports users to conduct their roles in communities (**user generated** data)

• ESFRI: European strategic forum of research infrastructure

- Upgrade of incoherent SCATter facility

EISCAT-3D



- Multidisciplinary seafloor observatory

EMSO



- Plate observing system

EPOS



- Global ocean observing infrastructure

EURO-ARGO



- Integrated carbon observation system

ICOS



- Biodiversity and ecosystem research infra

LIFEWATCH



• Diverse

standardises

- Terminologies
- Data models
- Metadata
- Service interfaces
- ...

• Semantic isolation

- Between metadata
- Between data content
- Between services
- ...

Metadata standards	Acquisition	Curation	Processing	Access	Community
SensorML	Y	Y		Y	
NetCDF			Y	Y	
ISO19115	Y	Y		Y	Y
ISO19156	Y	Y			
CSR	Y	Y			
Dublin Core				Y	
CERIF	Y	Y		Y	
CSMD	Y	Y		Y	Y
INSPIRE	Y	Y		Y	Y

ENVRI project

- Cluster project for environmental ESFRIs
- Identify common operations and needs among ESFRIs
- Guide the development of ESFRIs
- Promote Interoperability
- Enable interdisciplinary scientists to access, process, study and correlate data from multiple domains for *system level* research.



- Identify common requirements and operations from RIs
→ ODP approach
- Guide the RI development → Reference model approach
- Semantic gaps between RIs → ontological framework (OEILM)



Identify common requirements and operations

• Approach

- Collect use cases and requirements from all ESFRIs
- Analyse the use cases using the ODP approach
- Define a minimal set which crosses most of the ESFRIs

Functions/operations in the Data Curation Sub-system

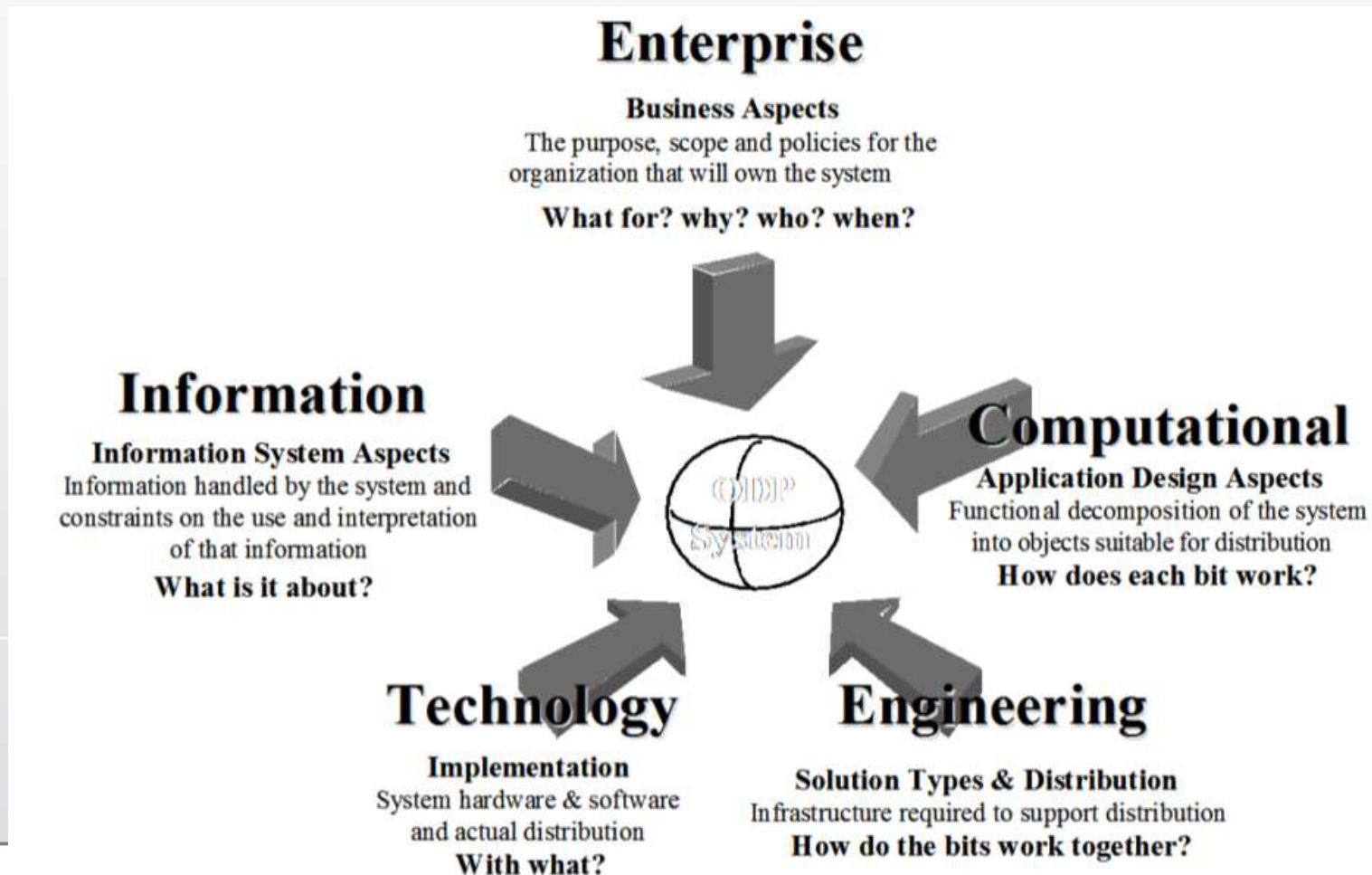
Functions/Embedded Services	ICOS	EPOS	EMSO	EISCAT-3D	LifeWatch	EURO-Argo
Data Quality Checking	Yes	Yes	Unknown	Yes	Not Applicable	Yes
Data Quality Verification	Yes	Unknown	Unknown	Unknown	Not Applicable	Yes
Data Identification	Yes	Yes	Yes	Unknown	Not Applicable	Unknown
Data Cataloguing	Unknown	Yes	Yes	Unknown	Not Applicable	Unknown
Data Product Generation	Yes	Yes	Yes	Yes	Not Applicable	Yes
Data Versioning	Yes	Unknown	Unknown	Unknown	Not Applicable	Unknown
Workflow Enactment	No	Yes	Unknown	Yes	Not Applicable	No
Data Preservation	Yes	Yes	Yes	Yes	Not Applicable	Yes
Data Replication	No	Yes	Unknown	Yes	Not Applicable	Yes
Data Replication Synchronisation	No	Unknown	No	Unknown	Not Applicable	Yes

Functions/operations at Data Access Sub system

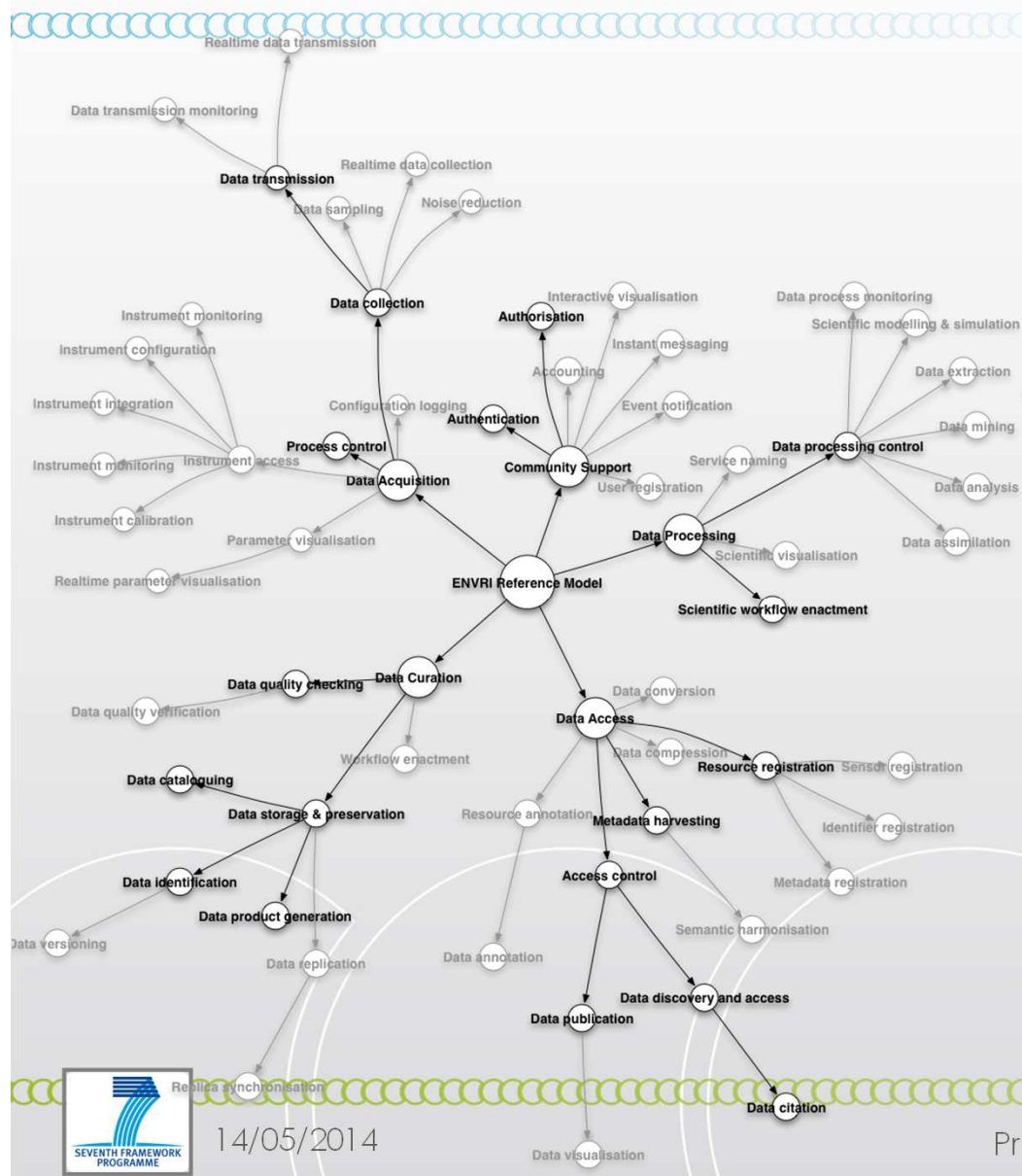
Functions/Embedded Services	ICOS	EPOS	EMSO	EISCAT-3D	LifeWatch	Euro-Argo
Access Control	Unknown	Yes	Unknown	Yes	Unknown	Unknown
Data Conversion	Yes	Yes	Yes	Yes	Yes	Yes
Data Compression	No	No	No	No	Yes	No
Data Visualisation	Yes	Yes	Yes	Yes	Yes	Yes
Data Publication	Yes	Unknown	Yes	Unknown	Yes	Yes
Data Citation	No	Unknown	Yes	No	Unknown	No
(Resources/Data) Annotation	Yes	Yes	Yes	No	Yes	Yes
Metadata Harvesting	Unknown	Unknown	Yes	No	Unknown	No
Resource Registration	Unknown	Yes	Yes	No	Yes	No
Semantic Harmonisation	No	Yes	Yes	No	Yes	No
Data Discovery and Access	Yes	Yes	Yes	Yes	Yes	Unknown

A full function list is on ENVRI wiki <http://envri.eu/group/envri/wiki/-/wiki/Main/Analyse%20Common%20Requirements%20for%20Data%20Processing>

- Open distributed processing (ODP): a multi viewpoint model for distributed systems. (ISO/IEC 10746)



Identify a minimal set



- Analysis of common requirements of ESFRI ENV RIs, resulted in **a set of common functionalities**
- Identified **a minimal model**
 - Focuses on **core interactions**
 - Represents the **most fundamental** functionalities
 - A **skeleton** which can be **extended**
 - Future development will be based on **community interests**

- **Derive use scenarios from common requirements, identifying *communities, roles, behaviours***
- **Model defines:**
 - **5 common *Communities*** in according to 5-subsystem
 - **Data Acquisition:** who collects raw data
 - **Data Curation:** who manages, archives quality data
 - **Data Publication:** who assists publication, discovery & access
 - **Data Service Provision:** who provides services to derive knowledge
 - **Data Usage:** who makes use of data/services
 - ***Community roles & behaviours***

- **Data-oriented approach:**

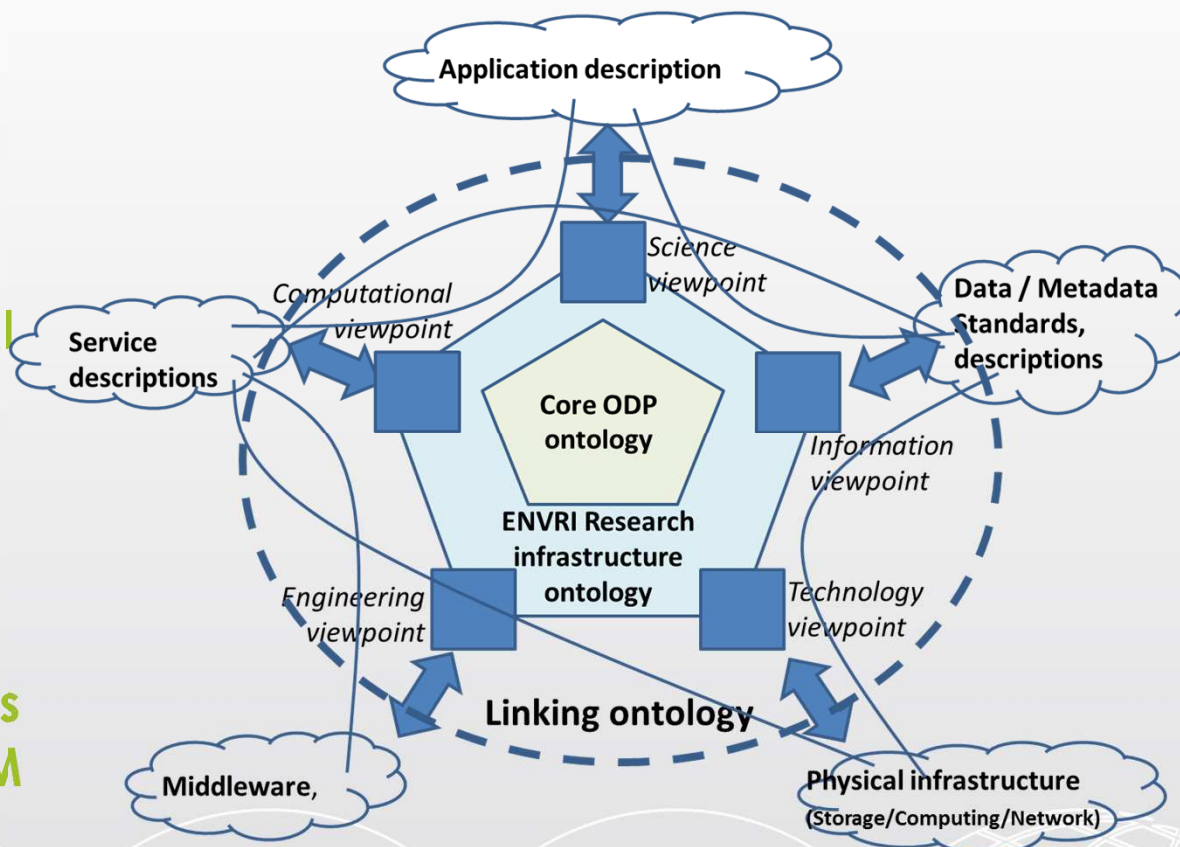
- Follow **data-lifecycle** in each subsystems
- Identify **information objects, actions, state changes** when events/actions occur

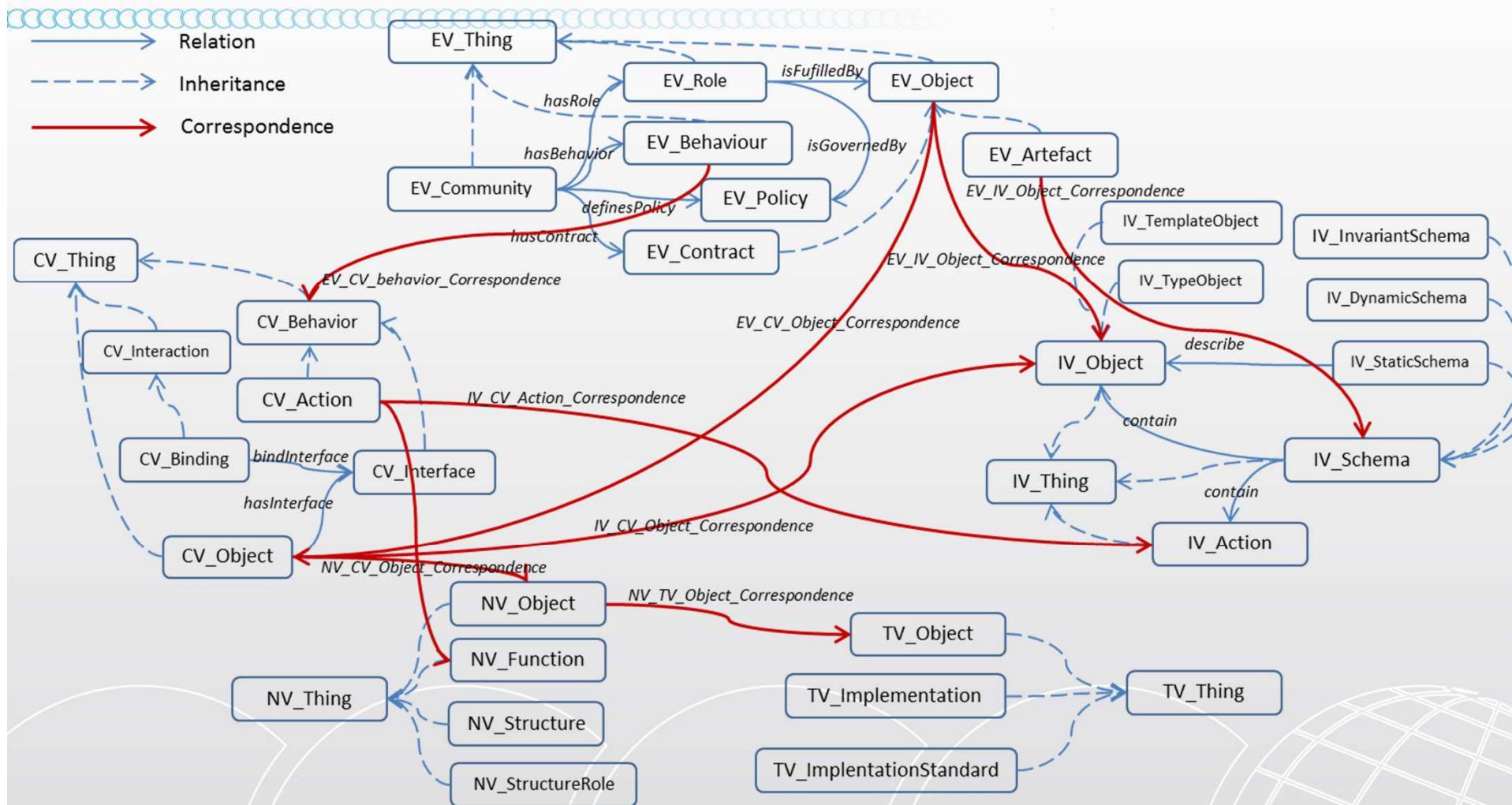
- **Model defines:**

- A set of **information objects** handled by a subsystem
- A set of **action types** that cause the states changes
- A set of **constraints** on those objects
- **Dynamic schema** -- how information objects evolve as the system operates
- **Static schema** -- allowable state changes

- **Service-oriented, Brokered approach**
 - **Core functionality is encapsulated** in a set of **service objects**
 - **Access** to such object **via brokers** which provides an **interoperability layer** between heterogeneous components
- **Model defines**
 - A set of **computational objects**
 - Each **encapsulates** specific functionalities
 - Each provides a set of **interfaces** to invoke functions
 - A set of **binding objects** to coordinate multi-party interactions

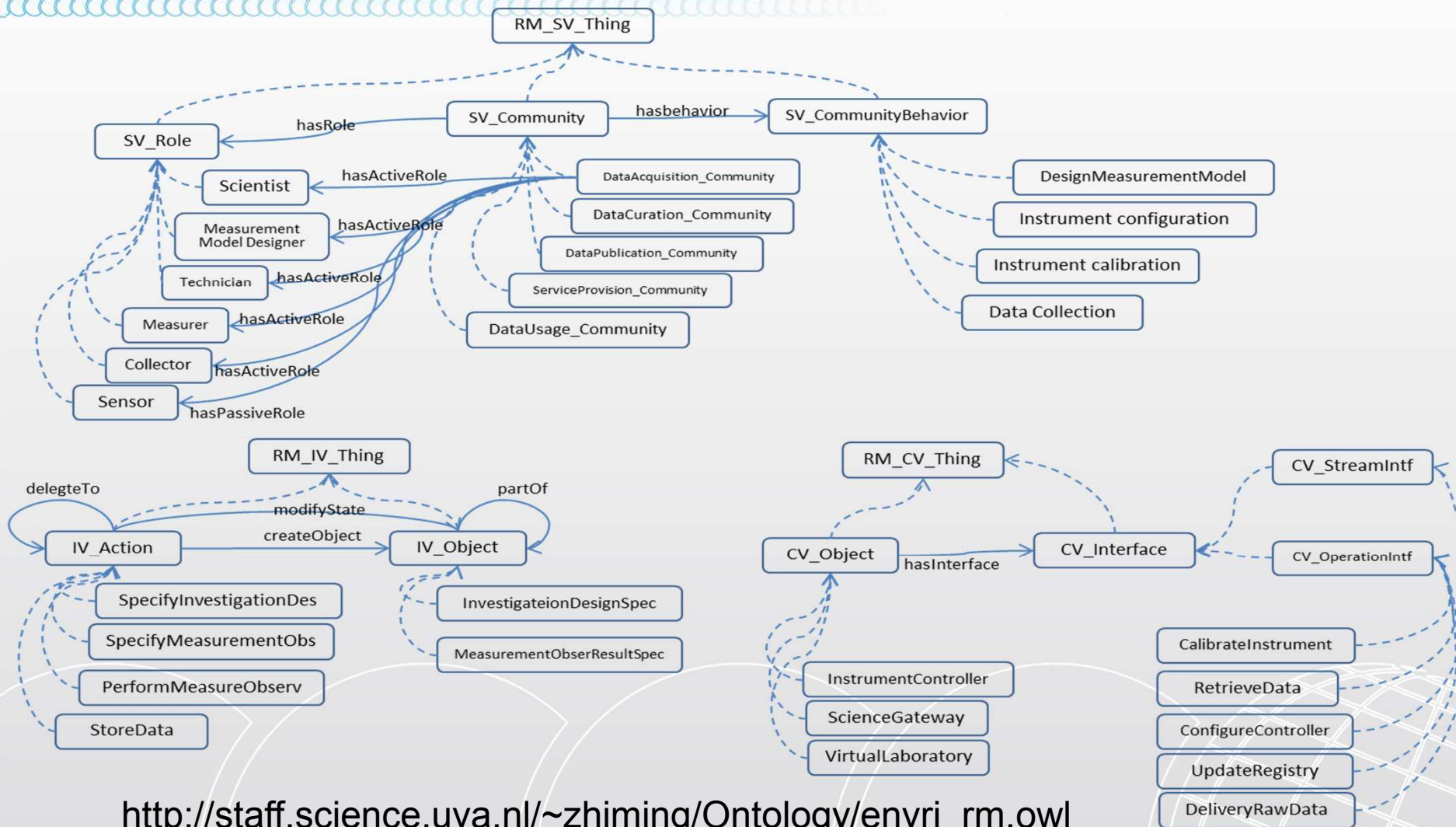
- **ODP ontology**
 - Vocabulary in ODP
 - Five viewpoints
- **ENVRI RM ontology**
 - Vocabulary in ENVRI RM
 - Five viewpoints extended from ODP
- **Linking ontology**
 - Concepts/properties extend the ENVRI RM from five viewpoints

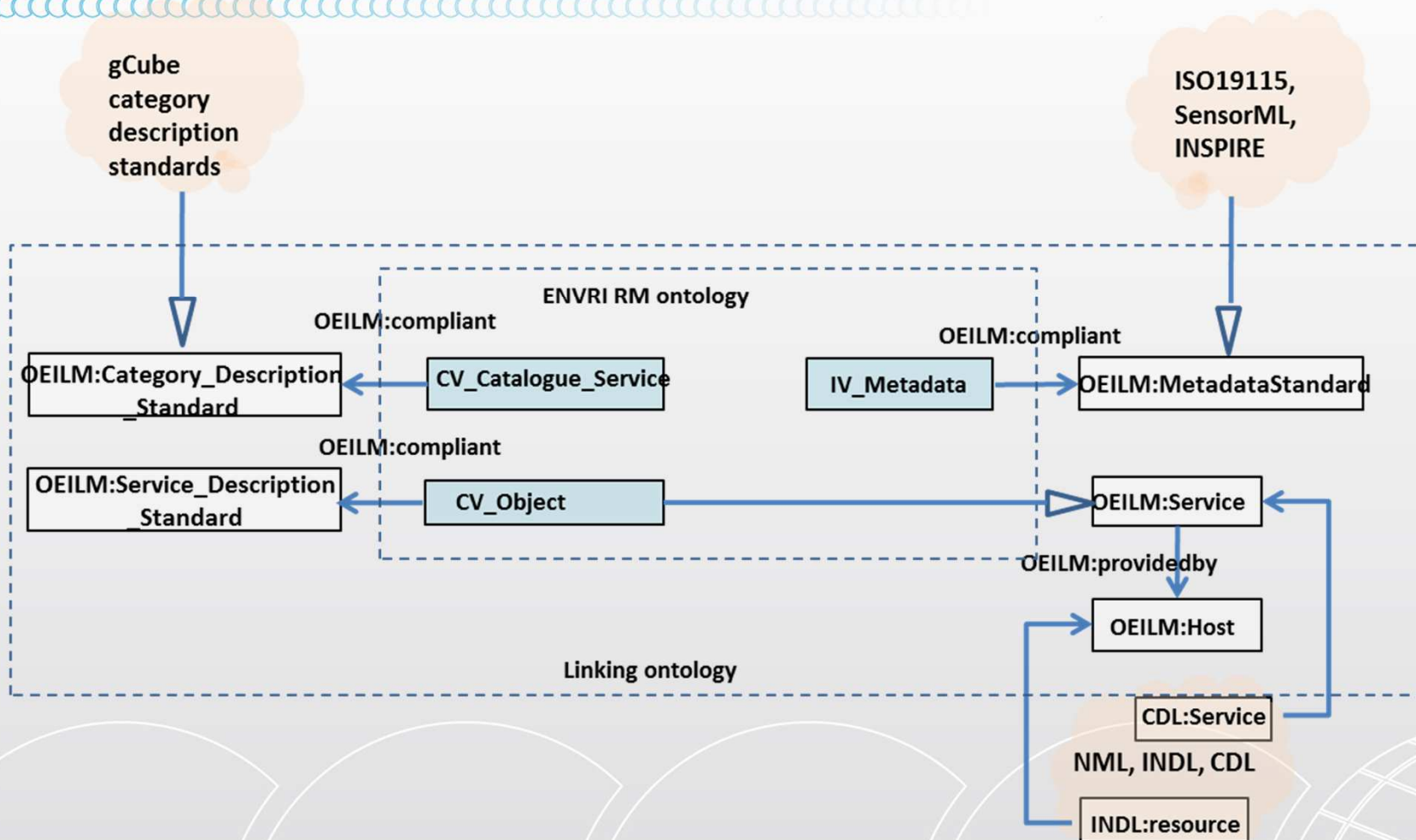




<http://staff.science.uva.nl/~zhiming/Ontology/odp.owl>

ENVRI RM ontology (part)







Inform

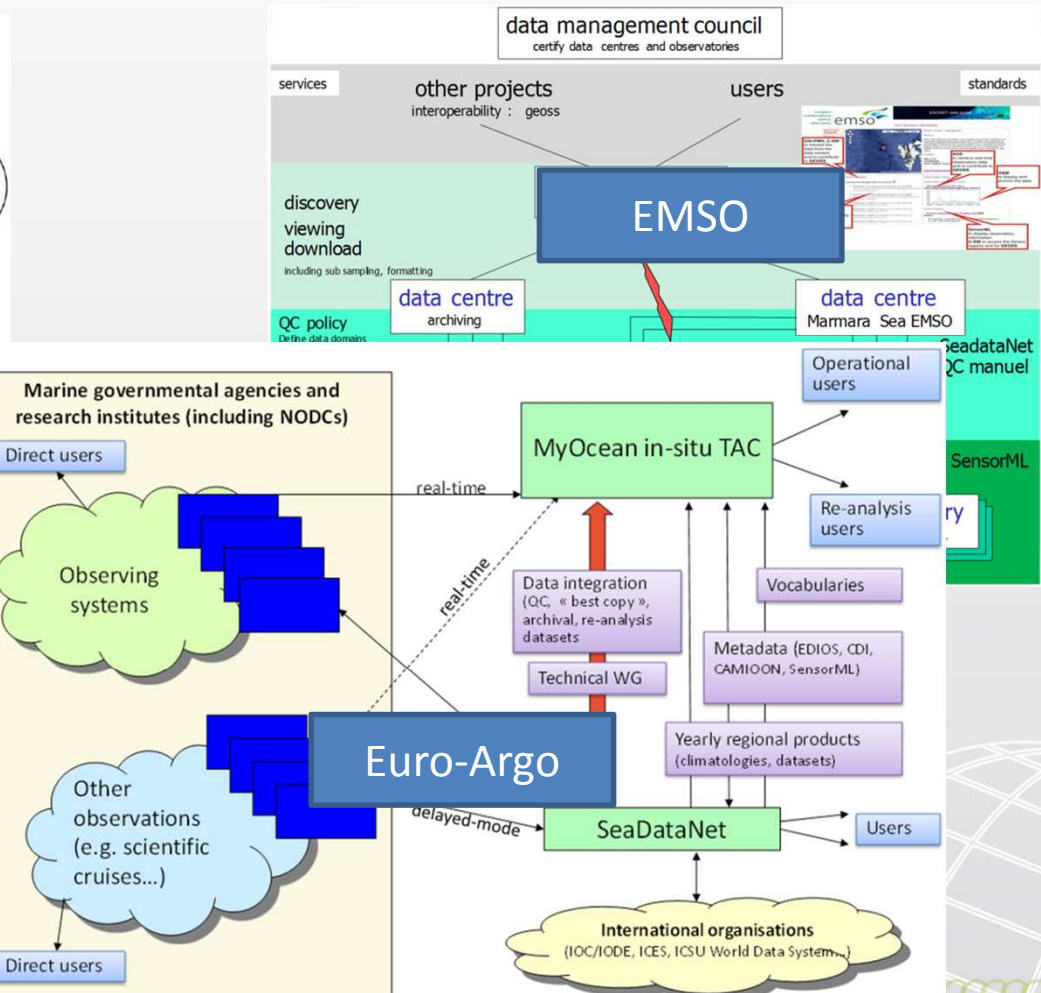
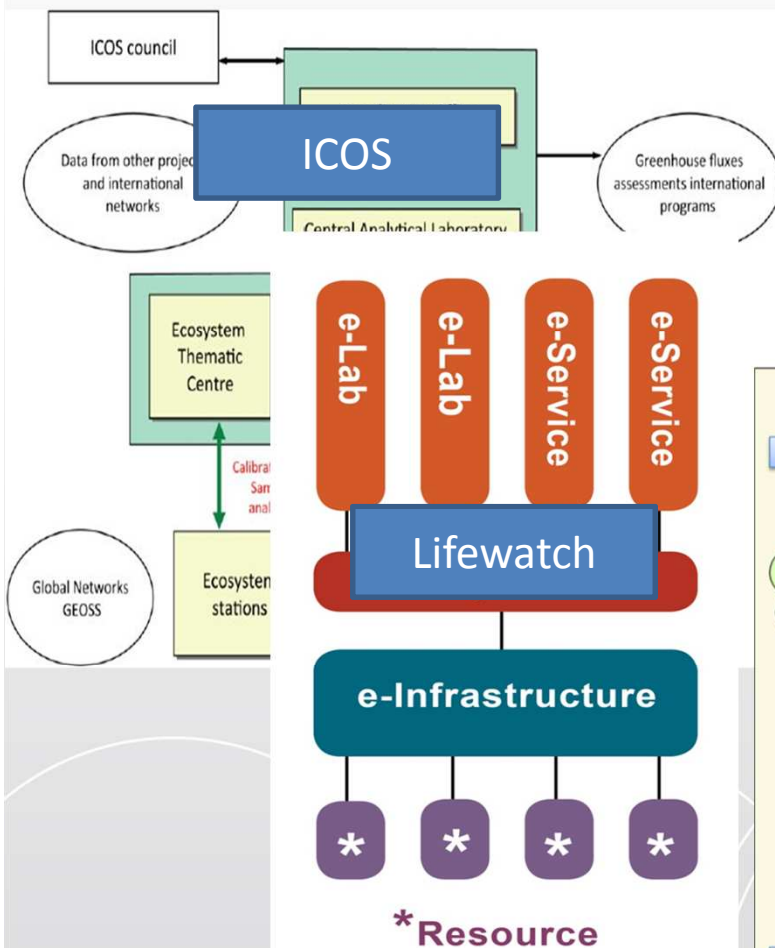


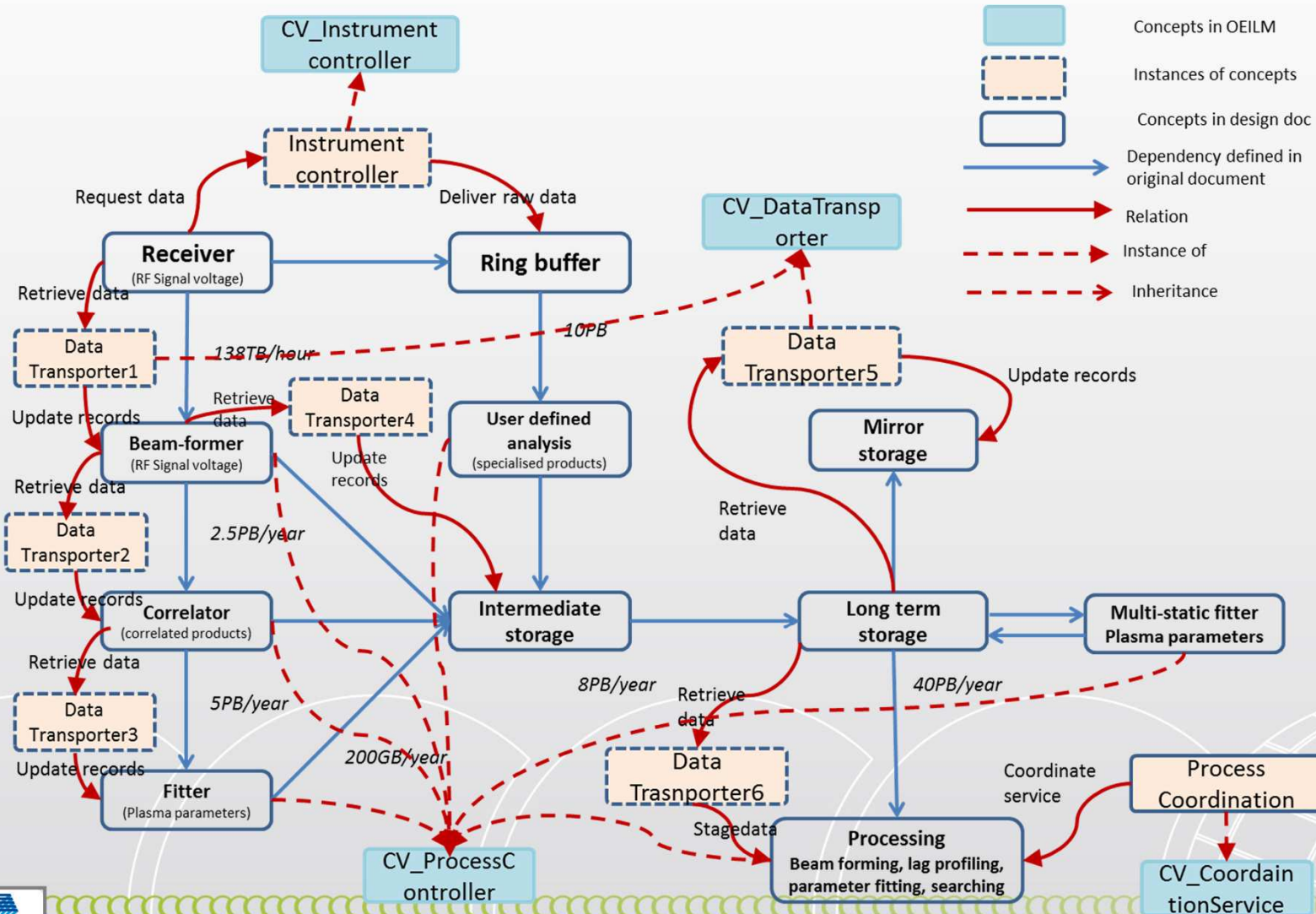
<http://envriontology.appspot.com/main/>.

- **ESFRI information sharing**
- **Resource discovery and workflow planning**

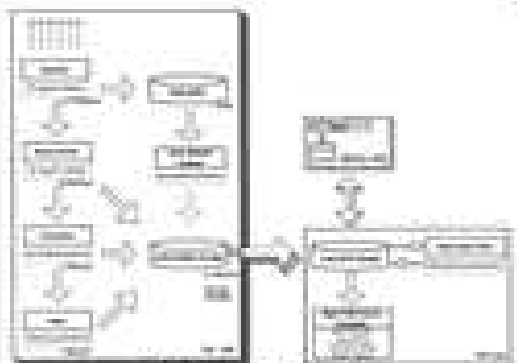
Use case 1: Sharing design documents between ESFRIs

• Different terminology makes sharing difficult





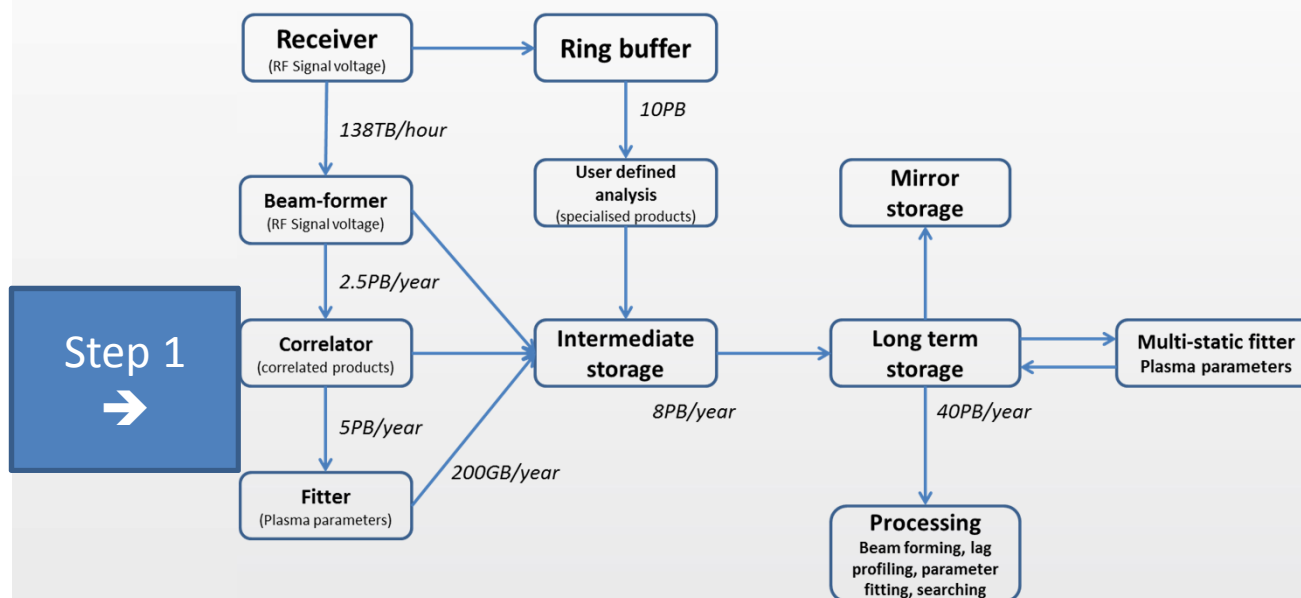
Step 1: design doc formalization



The architecture of ENVRI_3D is designed to support a wide range of scientific applications. It is a distributed system that can be scaled to meet the needs of different research groups. The system is designed to be flexible and adaptable, allowing for the integration of new components and the modification of existing ones. The system is designed to be robust and reliable, ensuring that data is not lost and that the system is always available to the users.

The system is designed to be scalable and flexible, allowing for the integration of new components and the modification of existing ones. The system is designed to be robust and reliable, ensuring that data is not lost and that the system is always available to the users. The system is designed to be scalable and flexible, allowing for the integration of new components and the modification of existing ones. The system is designed to be robust and reliable, ensuring that data is not lost and that the system is always available to the users.

The system is designed to be scalable and flexible, allowing for the integration of new components and the modification of existing ones. The system is designed to be robust and reliable, ensuring that data is not lost and that the system is always available to the users. The system is designed to be scalable and flexible, allowing for the integration of new components and the modification of existing ones. The system is designed to be robust and reliable, ensuring that data is not lost and that the system is always available to the users.

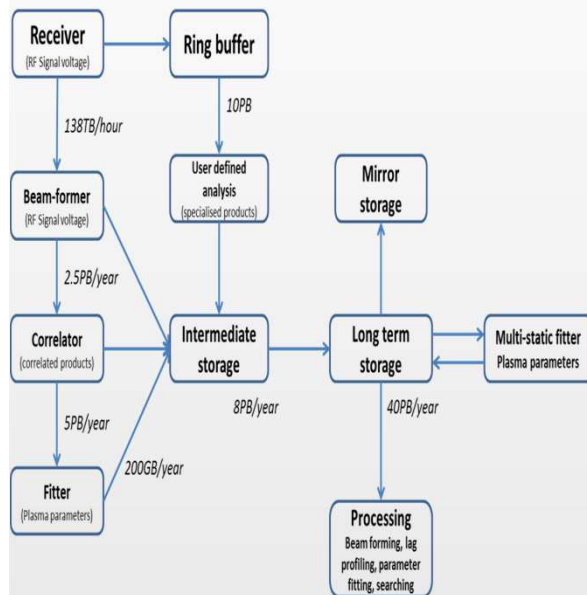


Step 1
→

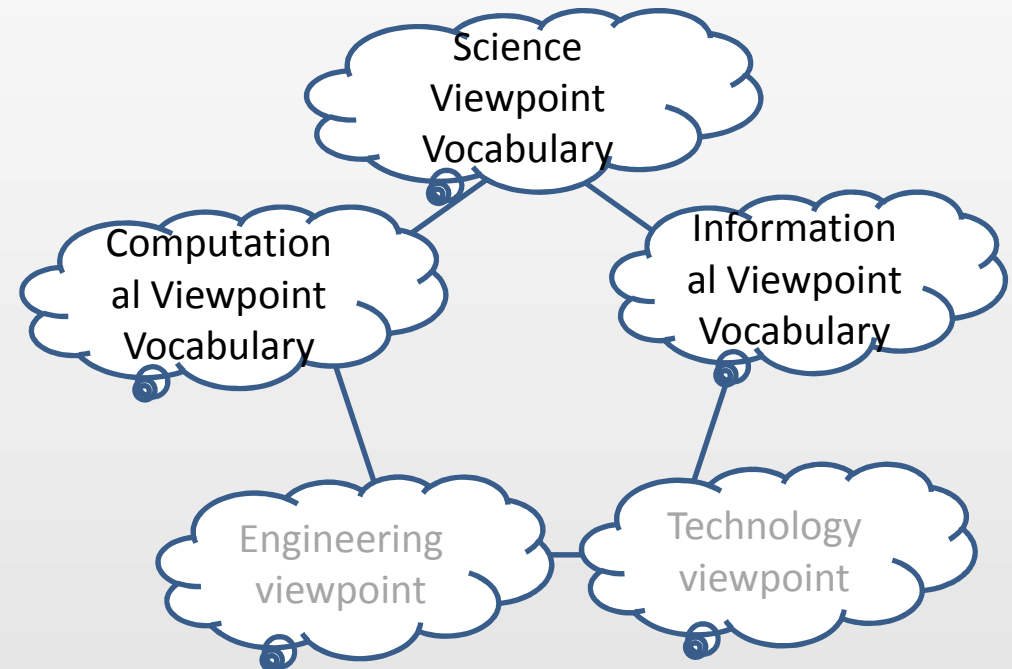
Identify the processes/data from the design documents

Provided by RI designers,
http://envri.eu/eiscat_3d-study-case

Step 2: find proper OEILM classes/properties (RM layer)



Step 2
→

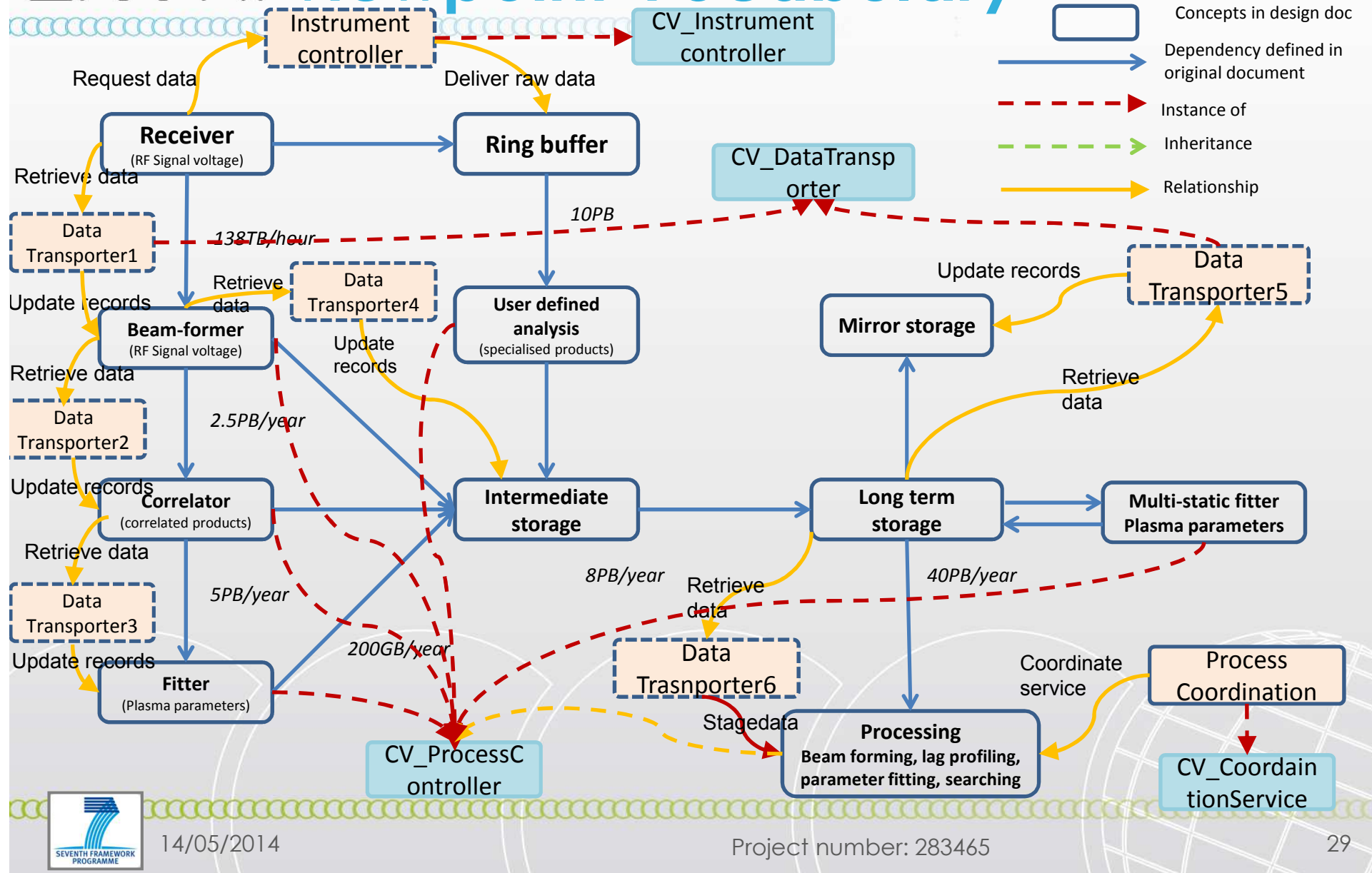


Map the concepts/data to the ENVRI RM concepts

- 1) Find the instances/classes
- 2) Identify the missing ones between instances

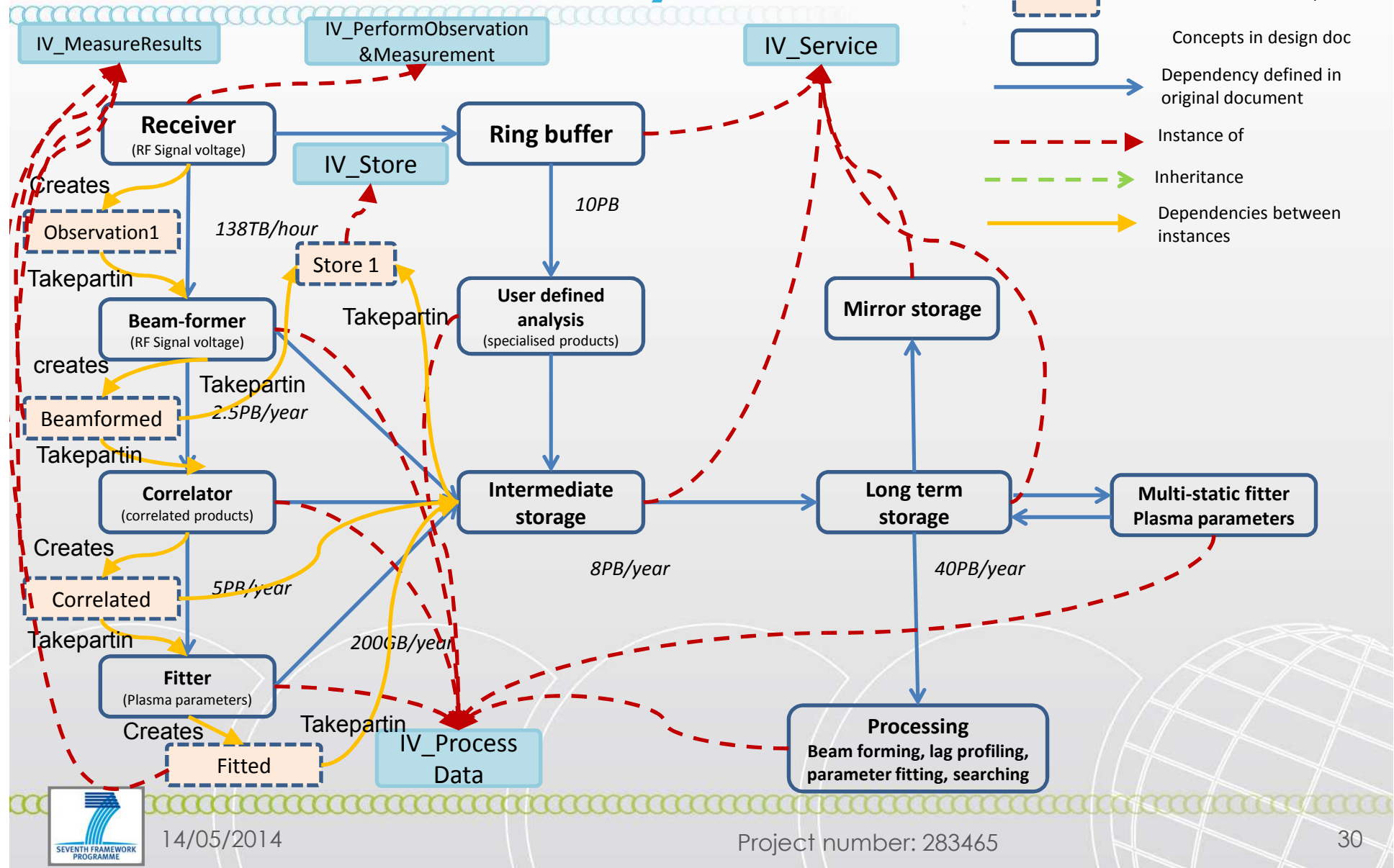


2.2 Using Computational viewpoint Vocabulary



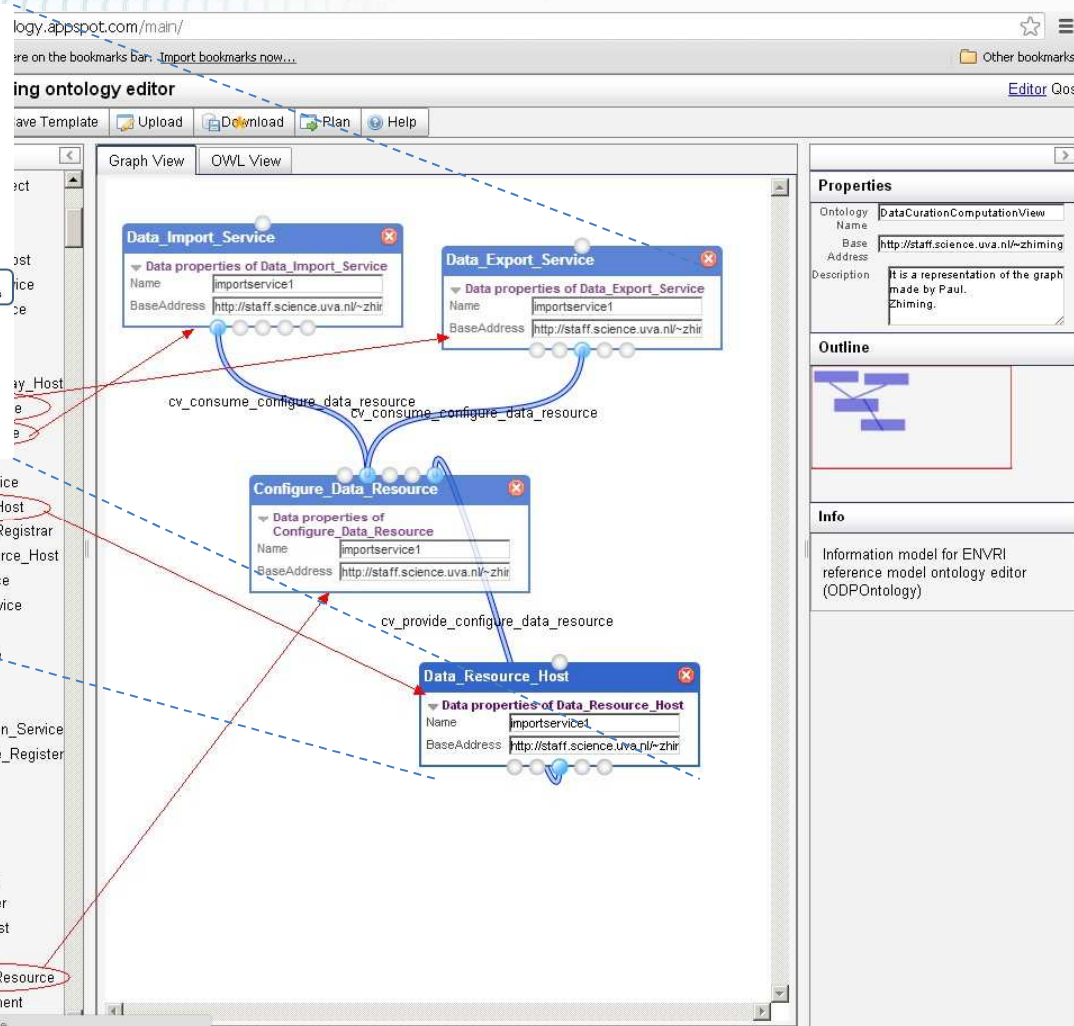
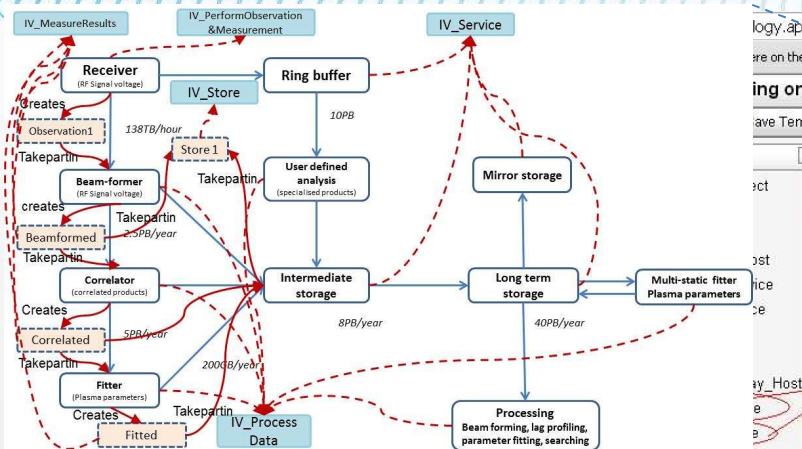
2.3 Using Information viewpoint

Vocabulary





OEILM tools for ENVRI annotation/description



14/05/2014

Project number: 283465

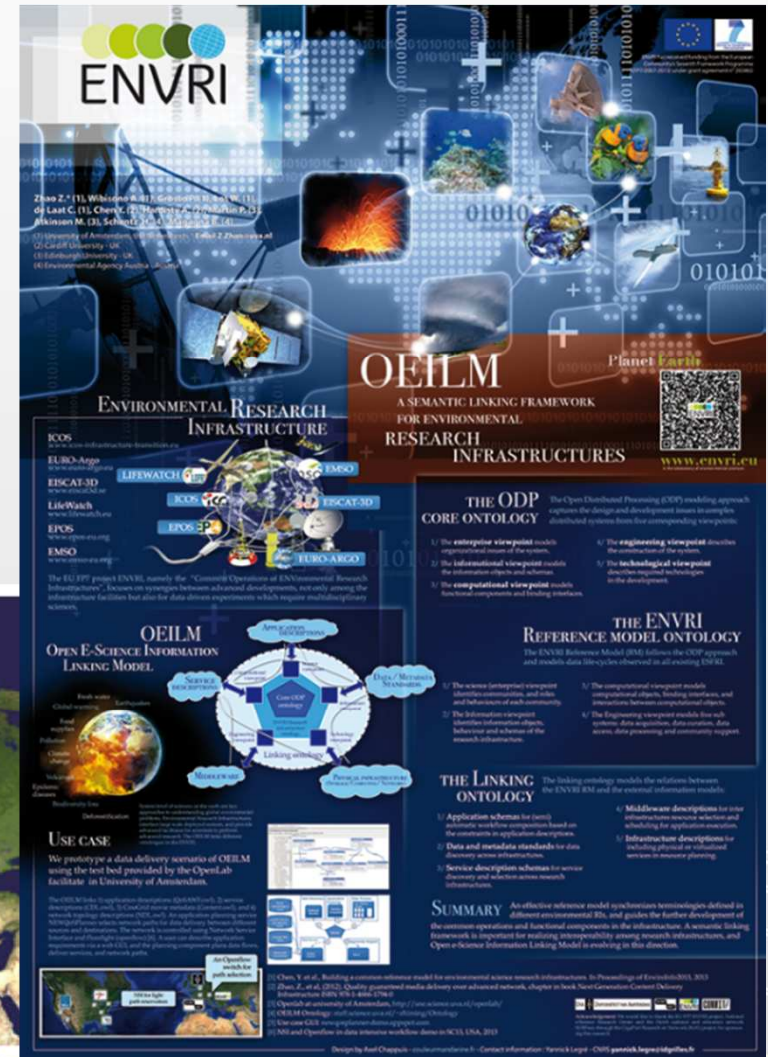
31

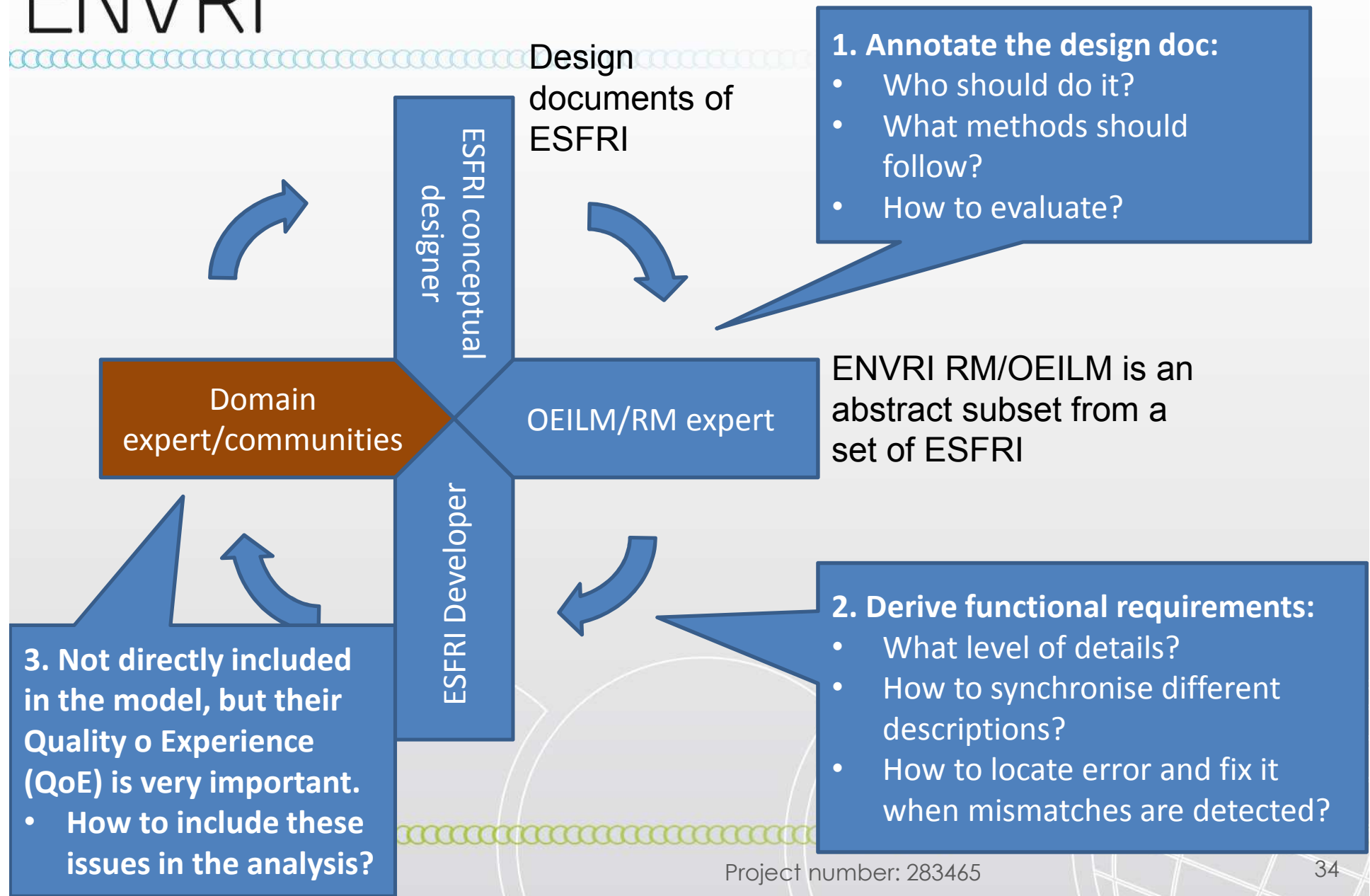
Step 3: Match-making

- From the annotation, we can identify the main functional components:
 - Science Viewpoint
 - SV_Role (SV_Sensor, SV_Storage)
 - SV_CommunityBehavior(SV_ProductGeneration, SV_ServiceComposition)
 - Computational Viewpoint
 - CV_InstrumentController, CV_DataTransporter, CV_ServiceCoordinator, CV_ProcessController
 - Information Viewpoint
 - IV_MeasurementResults, IV_Store, IV_Service, IV_ProcessData
- Functional requirements:
 - Acquisition:
 - SV_Sensor- CV_InstrumentController- IV_MeasurementResults
 - Curation
 - SV_Storage, CV_Transporter, *IV_Annotation, IV_PersistentIdentifier, IV_Metadata*
 - Access
 - SV_Storage, CV_Transporter, *IV_Metadata, IV_Identifier,*
 - Processing
 - SV_ProductionGeneration, SV_ServiceComposition, CV_Staging, CV_ServiceCoordinator, CV_ProcessController, IV_ProcessData, IV_Service,
 - Community
- EGI infrastructure
 - Curation (annotation, replication, catalogue)
 - Processing (computing power, parameter sweeping tool, vis, etc.)
 - Access (metadata, catalogue, AAA, storage, transfer)

ENVRI Use case 2: application planning using linking ontology

- Application QoS \leftrightarrow Services /Data \leftrightarrow Devices \leftrightarrow Network paths
- Live demo in SC 13





- The ENRI project will finish by the end of 2014
- The first version of OEILM has been made available for ENVRI ESFRI
- Ongoing work
 - Use cases with each specific ESFRI
 - Transfer knowledge to ESFRI
 - Collect feedback and refine OEILM for its version 2.
- Dissemination
 - RDA/EUDAT
- Exploitation
 - ENVRI 2 proposal

- ENVRI: www.envri.eu
- OEILM: <http://staff.science.uva.nl/~zhiming/Ontology>
- Zhao Z., Gross, P., Wouter. L., Chen Y., Hardisty, A., Martine, P., Magana, B., Schentz, H., (2013) *OEILM: a semantic linking framework for environmental research infrastructures*, Poster at IEEE e-Science 2013.
- Jiang, W., Zhao, Z., Grosso, P., de Laat, C., (2013) *Dynamic workflow planning on programmable infrastructure*, IEEE Int'l Conf. on Network Architecture Storage, 2013.
- Chen, Y., Martin, P., Schentz, H., Magagna, B., Zhao, Z., Hardisty, A., Preece, A., Atkinson, M., Huber, R., and Legre. R., *Building a common reference model for environmental science research infrastructures*. In Proceedings of EnviroInfo2013, 2013.
- Zhao, Z., van der Ham, J., Taal, A., Koning, R., Dumitru, C., Wibisono, A., Grosso, P., de Laat, C. (2012). *Planning data intensive workflows on inter-domain resources using the Network Service Interface (NSI)*, In the 7th Workshop on Workflows in Support of Large-Scale Science. in the context of Supercomputing 2012, Salt Lake City.
- Zhao, Z., Grosso, P. & Laat, C. de (2012). *OEIRM: An Open Distributed Processing based Interoperability Reference Model for e-Science*, Cloud&Grid interoperability workshop, Gwangju, Korean.
- Zhao, Z., Dumitru, C., Grosso, P. & Laat, C. de (2012). *Network resource control for data intensive applications in heterogeneous infrastructures*. the International Workshop on High Performance Data Intensive Computing, in 26th IEEE International Parallel and Distributed Processing Symposium, Shanghai.